1 [Title Page]

2 **Research on Multi-factor Forest Fire Prediction Model Using Machine Learning Method in**

3 **China**

4 **Yudong Li,[1] Zhongke Feng,[1] Ziyu Zhao,[1]Shilin Chen,[1]Hanyue Zhang,[1]**

5 **1 Precision Forestry Key Laboratory of Beijing, Beijing Forestry University, Beijing 100083,**

6 **China.**

7 **E-mail address ：**

8 **Yudong Li：lyd85842@163.com; Zhongke Feng：zhongkefeng@bjfu.edu.cn**

9 **Ziyu Zhao：zhaozy0315@126.com; Shilin Chen: chenshilin@bjfu.edu.cn**

10 **Hanyue Zhang：hanyue.zhang@foxmail.com**

11 **Correspondence should be addressed to Zhongke Feng: zhongkefeng@bjfu.edu.cn**

12 **Permanent address: Beijing Forestry University**

13 [Manuscript]

# 14 Research on Multi-factor Forest Fire Prediction Model

# 15 Using Machine Learning Method in China

16 Yudong Li,[1] Zhongke Feng,[1] Ziyu Zhao,[1]Shilin Chen,[1]Hanyue Zhang,[1]

17 1 Precision Forestry Key Laboratory of Beijing, Beijing Forestry University, Beijing 100083, China.

18 Correspondence should be addressed to Zhongke Feng: zhongkefeng@bjfu.edu.cn

## 19 Abstract

20 Forest fires can cause serious harm in many ways. Studying the scientific prediction of forest fires is an

21  important basis for preventing such fires. At present, there is little research on the prediction of long

22  time series forest fires in China. Choosing a suitable forest fire prediction model is of great importance

23  to China's forest fire prevention and control work. Based on data on fire hotspots, meteorology, terrain,

24  vegetation, infrastructure, and socio-economics collected from 2003 to 2016, we used a random forest

25  model as a feature-selection method to determine 13 major drivers of forest fires in China (such as

26  temperature, terrain etc.). The forest fire prediction models developed in this study are based on four

27  machine-learning algorithms: an artificial neural network, a radial basis function network, a

28  support-vector machine, and a random forest. The models were evaluated using the five performance

29  indicators of accuracy, precision, recall, f1 value, and area-under-the-curve value. We used the optimal

30  model to obtain the probability of forest fire occurrence in various provinces in China and create a

31  spatial distribution map of the areas with high incidences of forest fires. The results show that the

32  prediction accuracy of the four forest fire prediction models is between 75.8% and 89.2%, and the

33  area-under-the-curve value is between 0.840 and 0.960. The random forest model has the highest

34  accuracy (89.2%) and area-under-the-curve value (0.96). It is used as the optimal model to predict the

35  probability of forest fire occurrence in China. The prediction results indicate that the areas with high

36  incidences of forest fires are mainly concentrated in northeastern China (Heilongjiang Province and

37  northern Inner Mongolia Autonomous Region), southeastern China (including Fujian Province and

38  Jiangxi Province) etc. In those areas at high risk of forest fires, the management departments can

39  improve the forest fire prevention and control by establishing watch towers and using other monitoring

40  equipment. This study not only helps in understanding the main drivers of forest fires in China, but it

41  also provides a reference for the selection of high-precision forest fire prediction models and provides a

42  scientific basis for China's forest fire prevention and control work.

43     **Keywords:** forest fire occurrence in China; feature selection; forest fire driving factors; machine

44     learning; prediction model; forest fire prevention and control

# 1. Introduction

46     Forest fires are one of the most dangerous natural disasters. They have become the focus of worldwide

47     attention due to their rapid spread, their low controllability, and the hazards they pose[1-2]. Forest fires

48     have varying degrees of impact on human health and safety, the ecological environment and resources,

49     and society and the economy[3]. Forest fire prevention has therefore become a key research topic in the

50     fields of forestry and ecology[4-6].

51     The most effective way to control forest fires is to detect them quickly. Detection is usually divided

52     into three categories: satellite monitoring, smoke detection, and local perception (such as data analysis).

53     Satellite monitoring is expensive, involves delays, and is not fully applicable to all locations[7] . Smoke

54     detection also requires expensive equipment and maintenance work. In contrast, data for forest fire

55     analysis are easy to obtain, and data analysis is less expensive than the first two methods[8]. By

56     establishing a forest fire prediction model, we can predict the probability of the occurrence of a forest

57     fire and then strictly manage the area where the fire is likely to occur. This can directly reduce the

58     occurrence of forest fires and the potential casualties and economic losses[9-10]. This method is

59     therefore of great significance in forest fire prediction and prevention[11].

60     Much research has been conducted on forest fire prediction models. Logistic regression models are the

61     most commonly used. They have the advantage of solving the classification problem[12-16]. In recent

62     years, geographically weighted regression models have also been used. Wang[17]used a geographically

63     weighted regression model to predict regional fires in Gansu, China (2017). Guo[18] used a

64   geographically weighted logistic regression model to determine the relationship between human-made

65   fires and the potential drivers of forest fires in northern China(2016). Such a method can provide a

66   reasonable explanation for spatial heterogeneity, but the regression analyses can only be performed on

67   continuous variables; the method lacks analysis of categorical variables. Liu [19] used an exponential

68   equation to predict the number of forest fires in China, but this model analyzed only meteorological

69   factors(2017). Similarly, many researchers have used generalized linear regression models for forest

70   fire prediction. Miao et al.[20] used the zero-inflated Poisson model to predict the frequency of forest

71   fires in Japan in 2000(2008). Mandallaz et al. [21] used the Poisson model to predict forest fires in

72   France, Italy, etc. (1997). Guo et al. [22] used ordinary least squares regression, zero-inflated negative

73   binomial regression and the zero-inflated negative binomial model to predict the number of forest fires

74   in the Greater Xing'an Mountains area of Heilongjiang Province, China, and demonstrated that the

75   zero-inflated negative binomial model has the best performance(2010).

76   The development of artificial intelligence has led researchers to focus on building a forest fire

77   prediction model using machine-learning algorithms[23-31]. Artificial neural networks consist of

78   neurons with adjustable connection weights. Compared with traditional multiple linear regression

79   models or parametric regression models, neural networks have better self-organization and

80   self-learning capabilities. They have been widely used in forest fire prediction[32-34]. For example,

81   Maeda et al.[35] used artificial neural networks and multi-temporal images from MODIS/Terra-Aqua

82   sensors to detect areas at high risk of forest fires in the Amazon region of Brazil(2009). The results

83   showed that the error is less than 1, and the predictions are accurate. Sakr et al. [36] predicted the

84   occurrence of forest fires in developing countries through two meteorological factors using artificial

85   neural networks (2011).

86    A radial basis function (RBF) neural network is a three-layer neural network. It is a special case of back

87    propagation neural network. At present, little research has used RBF neural networks for forest fire

88    prediction. Samaher[25] used an RBF neural network to predict the forest fire risk in Portuguese

89    natural parks (2018). The final root mean square error was 54.2.

90    Support-vector machines (SVM) are most suitable for binary classification of data in the form of

91    supervised learning. They apply the principle of structural risk minimization and have good learning

92    ability. In recent years, researchers have begun to use SVMs to predict forest fires[8][37-40]. Samaher

93    [25] used five different soft computing (SC) technologies, including an SVM algorithm, to predict

94    areas at risk of forest fires (2018). He determined that the SVM algorithm provides more accurate

95    prediction than the other four algorithms. Cortez et al. [8] used five different Data Mining(DM)

96    algorithms to predict the area at risk of forest fires in the northeastern region of Portugal (2007). Their

97    results showed that the root mean square error was 64.7. Based on Cortez's research, Zhiqing et

98    al.[7]used the semi-definite programming model to select the optimal kernel function of the SVM to

99    establish an SVM model for forest fire prediction (2012). The mean square error was 1.76, and the

100   model effect was good.

101   The random forest (RF) algorithm is a well-known integrated learning algorithm. It can provide higher

102   accuracy than other algorithms. At present, the use of RFs to predict forest fires is relatively

103   established[41-43]. Liang et al. [44] used an RF model to predict the occurrence of forest fires in Fujian

104   Province, China, with an accuracy rate of 85% (2016). Pourtaghi et al. [45] used an RF algorithm to

105   study the sensitivity of forest fires in Golestan Province, Iran. Their results showed that the model

106   achieves the desired accuracy (2016).

107     Most of the current research focuses on certain areas, and there are few studies on the prediction and

108     analysis of long time series in China. Most such studies concentrate on the temporal and spatial

109     changes and influencing factors of forest fires in specific years[46-50]. The results of previous research

110     are therefore localized and limited, and there is a lack of research into the most suitable and

111     high-precision forest fire prediction model on the national scale.

112     In this study, we selected a variety of forest fire driving factors to build four prediction models based

113     on machine-learning algorithms. The models were evaluated using data on Chinese forest fires from

114     2003 to 2016. The study has three objectives: (1) identify the main forest fire driving factors and their

115     impact in China; (2) select the most suitable model for forest fire prediction in China by creating four

116     models and comparing and analyzing the fitting results; and (3) use the model that offers the most

117     accurate predictions to create a probability map of forest fires in China and put forward

118     recommendations for forest fire prevention.

## 119     2. Materials and Methods

### 120     2.1 Study Area and Data Resources

121     Located in east Asia on the west coast of the Pacific Ocean, China's territory is vast, with a total land

122     area of about 9.6 million square kilometers. The topography is high in the west, with vast mountains

123     and plateaus, and low in the east. The distance between the east and the west of the country is about

124     5,000 kilometers; the coastline of the mainland is more than 18,000 kilometers in length; and the

125     temperature and precipitation are diverse, forming a variety of climates[51]. The distribution of forest

126     resources in China is uneven, being mainly distributed in the northeast, south, and southwest regions.

127     The forested area is 220 million hectares, and the forest coverage rate is 22.96%.

128    The research data were divided into six parts: fire ignition data, meteorological data, terrain data,

129    vegetation data, infrastructure data, and socio-economic data. The fire point data were derived from

130    NASA's *Global Fire Atlas with Characteristics of Individual Fires, 2003–2016* (https://daac.ornl.gov/).

131    The *Global Fire Atlas* is a global dataset that tracks the daily dynamics of single fires. For each

132    individual fire, the dataset provides information about the fire's timing and location, scale, perimeter,

133    duration, speed, and direction of spread. These individual fire characteristics are based on the *Global*

134    *Fire Atlas* algorithms and estimated combustion day information from a 500-meter resolution product

135    of the 6 MCD64A1 combustion zone product of the Medium Resolution Imaging Spectroradiometer

136    (MODIS) collection.

137    This study used the fire point data for forest land in China from 2003 to 2016. The final number of fire

138    point data are 32746 (except Taiwan). The meteorological data are derived from the 14-day daily value

139    dataset of the China Meteorological Data Network (http://data.cma.cn/dataService/). The dataset

140    includes eight elements, such as the barometric pressure, temperature, relative humidity, and

141    precipitation of the station. Digital elevation model data were obtained through the Geospatial Data

142    Cloud website (http://www.gscloud.cn/). Vegetation data were represented by the normalized difference

143    vegetation index (NDVI), and the spatial distribution dataset of China's Quarterly Vegetation Index

144    comes from the Resource and Environment Data Cloud Platform (http://www.resdc.cn/). The basic

145    geographic data were taken from the "National Basic Geographic Database of 1:1 Million" on the

146    website of the National Geographic Information Resources Directory System. The data include the

147    locations of railways, highways, water systems, and residential areas. The socio-economic data include

148    population density and GDP per capita, and the grid data of the spatial distribution of population and

149    GDP were obtained from the Resource and Environment Data Cloud Platform. Figure 1 shows the map
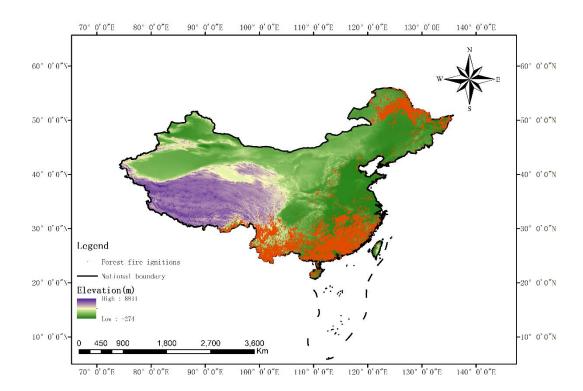
150    of the study area.

## 2.2 Data Preprocessing

### 2.2.1 Variable Handling

155    The dependent variable is a binary variable (i.e., whether a forest fire occurs), and so we used ArcGIS

156    to create a certain percentage of random points (non-fire points) and assigned 1 to fire points and 0 to

157    non-fire points[52]. To ensure that the data were not over-dispersed, random points were selected

158    according to experience in a ratio of 1:1[53], and, in principle, randomness in space and time should be

159    followed[54]. We used the ArcGIS 10.4 software to create random points and then used the 2015

160    National Land Use data as a basis to exclude random points that fell in bodies of water or urban land.

161    We obtained a total of 65,492 fire points and random points.

162    For the meteorological data, we first used the ArcGIS 10.4 software to match the sample points with

163    the nearest meteorological station through the Thiessen polygon method. We then extracted the

164    corresponding sample point weather data and used a SQL Server database to match the daily weather

165    data. For the terrain data, we used the spatial analysis tool in the ArcGIS 10.4 software to extract the

166    slope and aspect of the obtained digital elevation model data. Seasonal climatic differences have an

167    impact on vegetation status, and so we divided the year into spring (March, April, May), summer (June,

168    July, August), autumn (September, October, November), and winter (December, January, February)

169    [55-56]. We used the extraction and analysis tools of the ArcGIS software to extract NDVI data for the

170    sample points on an annual and quarterly basis.

171    Similarly, from the infrastructure data and socio-economic data, we extracted the information

172    corresponding to the sample points. We set the aspect and special festivals as categorical variables, and

173    the others as continuous variables. Table 1 shows the classification of aspect [57]. During certain

174    traditional festivals in China, people burn paper to commemorate their loved ones, which raises the

175    probability of a forest fire. We included as special festivals (value 1) the following dates: Chinese New

176    Year's Eve, the first day of the first lunar month, the second day of the first lunar month, the fifteenth

177    day of the first lunar month, and Qingming Festival and Zhongyuan Festival (July 15th of the lunar

178    calendar). Non-special festivals were set to 0.

179    **Table 1: Descriptions of aspect classifications**

| Aspect | Azimuth (degree) | Classification |
| --- | --- | --- |
| Gentle slope | -1 | 0 |
| Shady slope | 0~67.5, 337.5~360 | 1 |

| | | | |
|---|---|---|---|
| Semi-shady slope | 67.5~112.5, 292.5~337.5 | | 2 |
| Sunny slope | 157.5~247.5 | | 3 |
| Semi-sunny slope | 112.5~157.5, 247.5~292.5 | | 4 |

180    After processing, we obtained 20 independent variables and their possible values (see Table 2). Finally,

181    we performed data cleaning on the sample points and the various types of data extracted to remove

182    abnormal samples from the original dataset (including some samples with missing data and samples

183    with observations that were significantly outside the normal range).

184                         **Table 2: Descriptions of independent variables**

| Category | Independent Variable | Symbol | Variable Type |
|---|---|---|---|
| Location | Longitude (°) | *Lon* | Continuous Variable |
| | Latitude (°) | *Lat* | Continuous Variable |
| Terrain | Altitude (m) | *Alt* | Continuous Variable |
| | Slope (°) | *Slo* | Continuous Variable |
| | Aspect | *Asp* | Categorical Variable |
| Meteorology | Average surface temperature (°C) | *Avst* | Continuous Variable |
| | Daily maximum surface temperature (°C) | *Mast* | Continuous Variable |
| | Cumulative precipitation at 20–20 (mm) | *Pre* | Continuous Variable |
| | Average relative humidity (%) | *Arh* | Continuous Variable |
| | Hours of sunshine (h) | *Suh* | Continuous Variable |
| | Average temperature (°C) | *Ate* | Continuous Variable |
| | Daily maximum temperature (°C) | *Mate* | Continuous Variable |
| | Average wind speed (m/s) | *Aws* | Continuous Variable |

| | Maximum wind speed (m/s) | *Mws* | Continuous Variable |
|---|---|---|---|
| Infrastructure | Distance from fire point to highway (m) | *Hig* | Continuous Variable |
| | Closest distance from fire point to residential area (m) | *Set* | Continuous Variable |
| Social humanity | Population | *Pop* | Continuous Variable |
| | GDP | *GDP* | Continuous Variable |
| | Special festival | *Sfe* | Categorical Variable |
| Vegetation | NDVI | *NDVI* | Continuous Variable |

185　**2.2.2 Data Normalization**

186　Given the different dimensions and magnitudes of the factors above, the data were normalized to

187　eliminate the variation in dimensions, avoid large differences in the magnitudes of the input and output

188　data, and balance the contributions of various factors. All the data were converted to between 0 and 1.

189　Table 3 shows the normalized formulas and specific interpretations of the independent variables.

190　**Table 3: Normalized formulas and explanations**

| No. | Formula | Explanation | Variables using this formula |
|---|---|---|---|
| (1) | $x_i{}^* = \dfrac{x_i - x_{min}}{x_{max} - x_{min}}$ | $x_i$ and $x_i{}^*$ are the values before and after data normalization; $x_{max}$ and $x_{min}$ are the maximum and minimum values of the full sample data. | *Lon, Lat, Alt, Avst, Mast, Pre, Suh, Ate, Mate, Aws, Mws, Hig, Set, Pop, GDP* |
| (2) | $x_\alpha = \sin\alpha$ | $\alpha$ is the slope value | *Slo* |
| (3) | $x_\gamma = \dfrac{\gamma}{100}$ | $\gamma$ is the humidity value | *Arh* |

## 2.3 Research Method

### 2.3.1 Artificial Neural Networks

Artificial neural networks (ANN) have become widely used in feedforward networks due to their clear structure, fast operation, easy implementation, and abilities for self-learning and adaption to the environment [24][58-59]. ANNs consist of three parts: an input layer, an output layer, and a hidden layer. The hidden layer may be a topological structure of one or more layers, as shown in Figure 2. The input layer does not perform any calculations. It is used to receive data; that is, to transfer data to the adjacent hidden layer with different weights. The hidden layer processes the data through a nonlinear activation function and then passes it to the output layer. The final result is obtained from the output layer [60]. The mathematical principle is as follows:

$$\begin{cases} h^{(1)} = \varphi^{(1)}(\sum_{i=1}^{n} x_i \cdot \omega_j^{(1)} + b^{(1)} \\ y = \varphi^{(2)}(\sum_{j=1}^{n} h_i^{(1)} \cdot \omega_j^{(2)} + b^{(2)} \end{cases} \tag{4}$$

In the formula, input layer $x \in R^m$, hidden layer output $h \in R^n$, output layer $y \in R^K$, input layer to hidden layer weight connection matrix $\omega^{(1)} \in R^{m \times n}$, the weight connection bias from the input layer to the hidden layer $b^{(1)} \in R^n$, the weight connection matrix and the bias from the hidden layer to the output layer are $\omega^{(2)} \in R^{n \times K}$ and $b^{(2)} \in R^{n \times K}$.
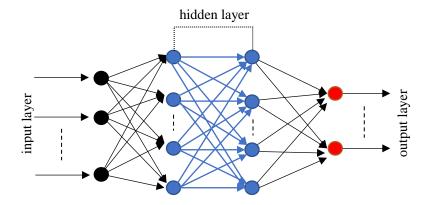
206

**Fig. 2 Diagram of the structure of an ANN**

207

**2.3.2 Radial Basis Function Neural Network**

208

The radial basis function (RBF) neural network structure is a feedforward structure with an input layer,

209

a single hidden layer, and an output layer [61]. Its advantages are concise training and fast learning

210

convergence speed, which can approximate any nonlinear function. It has been widely used in

211

time-series forecasting, nonlinear control systems, and the graphics-processing field. The basic idea of

212

an RBF network is as follows. The RBF is used as the "base" of the hidden unit to form the hidden

213

layer space. The hidden layer transforms the input vector and transforms the low-dimensional pattern

214

input data into the high-dimensional space. The result is that the data are linearly separable in the

215

high-dimensional space. The output of the RBF neural network is:

216

$$y_i = \sum_{i=1}^{h} \omega_{ij} \exp\left(-\frac{1}{2\sigma^2}\left\|x_p - c_i\right\|^2\right) \quad j = 1,2,\cdots,n \tag{5}$$

217

where $x_p = (x_1{}^p, x_2{}^p, \cdots, x_m{}^p)^T$ is the $p$-th input sample (p = 1,2,3,...,P), $P$ is the total number of

218

samples, $c_i$ is the center of the hidden layer node of the network, $\omega_{ij}$ is the connection weight from

219

the hidden layer to the output layer, $i = 1,2,3,...,$h is the number of hidden layer nodes, $y_i$ is the actual

220

output of the $j$-th output node of the network corresponding to the input sample [62].

221

222    ### 2.3.3 Support-Vector Machines

223    Support-vector machines (SVM) are mainly used for pattern classification and nonlinear regression.

224    They are general learning algorithms based on the principle of structural risk minimization. The core

225    idea of SVMs is to establish a classification hyperplane as a decision surface to maximize the isolation

226    edge between the positive and negative examples, thereby providing a high generalization

227    performance[63]. SVMs can improve the ability to transform data from high-dimensional spaces by

228    flexibly using kernel functions when dealing with various nonlinear problems. Taking a two-class SVM

229    as an example, given a training set $T = \{(x_1, y_1), \cdots (x_l, y_l)\} \in (X \times Y)^l$, where $x_i \in X = R^n, y_i \in$

230    $\{1, -1\}(i = 1, 2, \cdots l)$, $x_i$ is the feature vector. The penalty parameter $C$ and the kernel function

231    $K(x, x')$ are first selected, and the optimization problem is then constructed and solved as follows [62]:

232
$$\min_\alpha \frac{1}{2} \sum_{i=1}^{j} \sum_{j=1}^{l} y_i y_j a_i a_j K(x, x') - \sum_{j=1}^{l} \alpha_j \tag{6}$$

233
$$s.t. \ \sum_{i=1}^{l} y_i \alpha_i = 0, \ 0 \leq \alpha_i \leq C, i = 1, \cdots, l \tag{7}$$

234    The optimal solution is then obtained: $\alpha^* = (\alpha_1^*, \cdots, \alpha_l^*)^T$. A positive component of $\alpha^*$ :$0 \leq \alpha_j^* \leq C$

235    is then selected, and the threshold is calculated as follows:

236
$$b^* = y_j - \sum_{i=1}^{l} y_i \alpha_i^{\cdot} K(x_i - x_j) \tag{8}$$

237    Finally, the decision function is constructed:

238
$$f(x) = sgn(\sum_{i=1}^{l} \alpha_i^* y_i K(x, x_i) + b^*) \tag{9}$$

239    ### 2.3.4 Random Forest

240    A random forest is a highly flexible machine-learning algorithm with wide application prospects. In

241 essence, an RF is a classifier consisting of multiple decision trees formed by random methods. These

242 trees are not related, hence its alternative name: "random decision tree." When the test data enter the

243 RF, each decision tree is classified, and the category with the most classification results among all the

244 decision trees is taken as the final result. The RF algorithm has the following advantages: it evaluates

245 the importance of each feature in classification problems, it can process input samples with

246 high-dimensional features, and it does not require a reduction in dimensionality. The method is as

247 follows.

248 Let $N$ be the number of attributes of the sample. $n$ is an integer greater than 0 and less than $N$. First, the

249 bootstrap method is used for resampling, randomly generating $M$ training sets S1, S2, ...SM. The

250 decision tree A1, A2, ...AM corresponding to each training set is then generated. Before selecting the

251 attribute in each non-leaf node, $n$ attributes are randomly selected from the $N$ attributes as the split

252 attribute set of the current node, and the node is split in the best split mode among the $n$ attributes. Each

253 tree grows intact without pruning. For the test set sample $X$, each decision tree is used to test and obtain

254 the corresponding categories C1(X), C2(X), ...CM(X). Finally, the voting method is adopted, and the

255 category with the most output among the $M$ decision trees is regarded as the category to which the test

256 set sample $X$ belongs [62].

257 **2.3.5 Model Performance Evaluation**

258 In this study, we used five performance indicators: accuracy, precision, recall, f1 value, and

259 area-under-the-curve (AUC) value to evaluate the performance of the models. Descriptions of the five

260 indicators are given below.

261 1. Accuracy: the proportion of the number of samples (TP and TN) that are correctly predicted to the

262    total number of samples. The formula is as follows:

$$P = \frac{TP+TN}{TP+FP+TN+FN} \tag{10}$$

264    2. Precision: characterizes the classification effect of the classifier, which is the correct frequency value

265    predicted in the instance of the positive sample:

$$T = \frac{TP}{TP+FP} \tag{11}$$

267    3. Recall: characterizes the recall effect of a certain class. It is the correct frequency of prediction in the

268    instance of the label as the positive sample:

$$R = \frac{TP}{TP+FN} \tag{12}$$

270    4. f1 value: the value used to measure precision and recall. It is the harmonic mean of these two values:

$$f1 = \frac{(1+a^2)PR}{a^2(P+R)} \tag{13}$$

272    5. A ROC (receiver operating characteristic) curve is a method to judge the prediction effect of the

273    model[63]. The prediction accuracy of the model is judged by the value of the area under the curve

274    (AUC). The AUC ranges from 0.5 to 1. The larger the value, the closer the fit of the model.

275    Note: TP, FN, FP, TN in the formulas are the labels of the confusion matrix form of the output result.

276    The form is shown in Table 4:

**Table 4:    Confusion matrix form**

| Prediction (column) / label (row) | Positive sample | Negative sample |
| --- | --- | --- |

| | | | |
|---|---|---|---|
| Positive sample | TP | FN |
| Negative sample | FP | TN |

# 3. Results

In this study, we used the MATLAB (MathWorks, USA, MATLAB 2019a) [58] and RStudio (JJ Allaire, RStudio-1.2.5042/R 3.6.3) programming languages to implement the algorithms. We used MATLAB to build the SVM, ANN, and RBFN models and used RStudio to build the RF models.

To evaluate feature factors and model performance issues, the dataset was divided into two parts by randomly selecting 70% of the preprocessed sample data as the training set and 30% as the test set [59].

## 3.1 Feature Selection

We used the RF algorithm to perform feature selection on all variables after preprocessing, and we selected the subset of features that have the greatest impact on the dependent variable for the next model-building process. We divided the full sample according to the above-mentioned proportions (70% of the training set and 30% of the test set), and we obtained five training samples after repeating the process five times. Then we used the varSelRF package in the R language to perform feature variable selection calculations on the five training samples to obtain the variable subsets of the five intermediate models. Finally, we chose the variables that appeared 3 times or more in the five variable subsets as the main forest fire driving factors to enter the model fitting process. Table 5 shows the results of feature selection.

**Table 5: Results of variable selection based on RF**

| No. | Variable | Sample 1 | Sample 2 | Sample 3 | Sample 4 | Sample 5 | Frequency |
|---|---|---|---|---|---|---|---|

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 1 | Lat | + | + | + | + | + | 5 |
| 2 | Lon | + | + | + | + | + | 5 |
| 3 | Avst | + | + | + | + | + | 5 |
| 4 | Mast | + | + | | + | + | 4 |
| 5 | Pre | + | + | | + | + | 4 |
| 6 | Arh | + | + | + | + | + | 5 |
| 7 | Suh | + | + | + | + | + | 5 |
| 8 | Ate | + | + | + | + | + | 5 |
| 9 | Mate | + | + | + | + | + | 5 |
| 10 | Aws | | | | | | 0 |
| 11 | Mws | | | | | | 0 |
| 12 | Alt | + | + | + | + | + | 5 |
| 13 | Slo | | | | | | 0 |
| 14 | Asp | | | | | | 0 |
| 15 | Set | | | | | | 0 |
| 16 | Hig | | | | | | 0 |
| 17 | GDP | + | + | | + | + | 5 |
| 18 | Pop | + | + | + | + | + | 5 |
| 19 | NDVI | + | + | + | + | + | 5 |
| 20 | Sfe | | | | | | 0 |

295    The results show that the main influencing variables are longitude, latitude, average surface

296    temperature, daily maximum surface temperature, accumulated precipitation, average relative humidity,

297    sunshine hours, average temperature, daily maximum temperature, altitude, population, GDP, and

298    NDVI. These variables performed subsequent model fitting. Then the mean decrease accuracy obtained

299    by the RF algorithm was used to evaluate the importance of the variable. The larger the value, the

300    greater the importance of the variable. Figure 3 shows the importance of each variable in five random

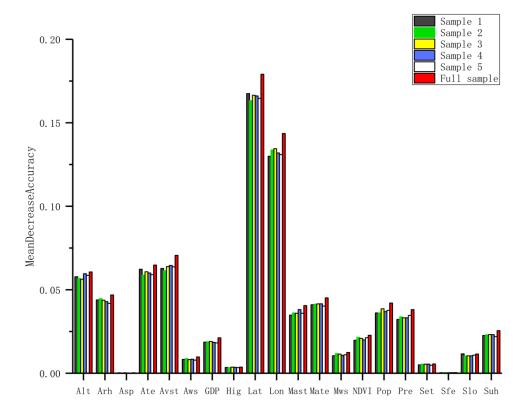301    training samples and 20 feature subsets in the full sample.



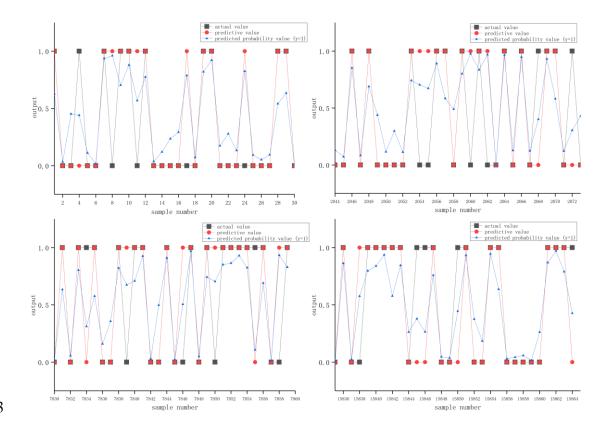303                              **Fig. 3 Feature subset importance**

304    It can be seen from Figure 3 that the longitude and latitude have the greatest influence on the

305    occurrence of forest fires. Thus, the location factor has the most influence on the occurrence of forest

306    fires. In addition, altitude also affects the occurrence of forest fires. The second greatest influence is the

307    temperature factor (average surface temperature and average temperature), reflecting the fact that high

308    temperatures can cause fires. Meteorological factors, including rainfall, sunshine hours, and average

309    relative humidity, can cause forest fires to varying degrees. Human activities (GDP and population) and

310    vegetation coverage also have an influence on the occurrence of forest fires but less so than other

311    factors such as weather and location. The variables not selected by the RF algorithm include average

312    wind speed, maximum wind speed, aspect, slope, the closest distance from the fire point to the highway,

313    the closest distance from the fire point to a residential area, and special festivals. The results indicate

314    that these seven variables have little influence on the occurrence of forest fires during the data analysis.

## 3.2 Model Fitting Results

### 3.2.1 Artificial Neural Network

317    The input layer of the ANN consists of 13 neurons after feature selection: Lat, Lon, Avst, Mast, Pre,

318    Arh, Suh, Ate, Mate, Alt, GDP, Pop, and NDVI. The output layer contains two units (1 or 0). We used

319    gradient descent to optimize the algorithm. Finally, we used a single hidden layer containing five units.

320    The comparison between the predictive value and the actual value in the test dataset is shown in Figure

321    4. Note: Due to the large sample size, only a part of the sample comparison chart is displayed. This is

322    also the case for the following comparison charts.



323

324                      **Fig. 4 Comparison charts of the predictive and actual values of the ANN**

325      **3.2.2 Radial Basis Function Neural Network**

326      The input and output layer variables of the RBF neural network were the same as those of the ANN.

327      After training, we obtained a hidden layer containing 10 units. The comparison charts of the predictive

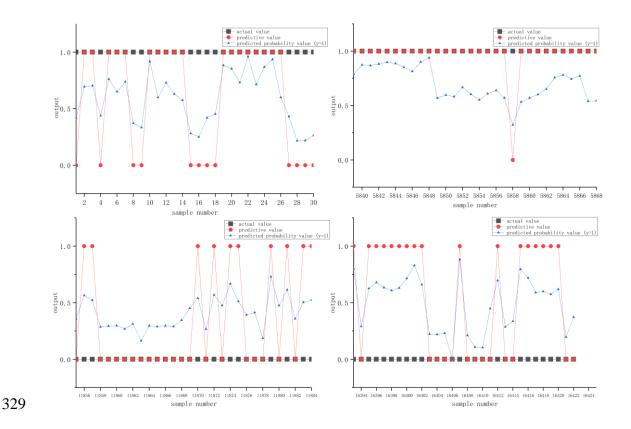328      and actual values of the test set are shown in Figure 5.

329



330                **Fig. 5 Comparison charts of the predictive and actual values of the RBFN (part of the**

331                              **sample)**

332      **3.2.3 Support-Vector Machine**

333      We used the LIBSVM package of the MATLAB software to construct the SVM. The model was

334      constructed using the RBF kernel function for processing nonlinear data. We used the grid search

335    method and 10-fold cross validation to select parameters and determine the penalty parameter $C$ and

336    the kernel parameter $g$. Figure 6 shows a contour map and a 3D view of the result of the SVC

337    parameter selection. After calculation, the accuracy rate of the grid search method reached 83.9%, and
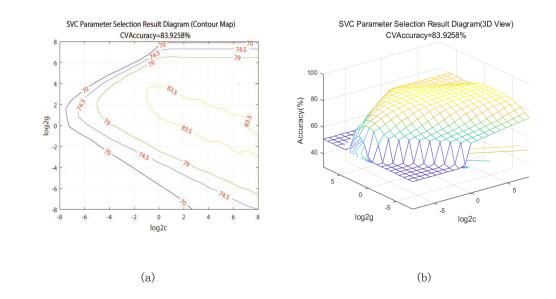
338    the accuracy rate of cross-validation reached 82.6%.

339



340                              (a)                                                              (b)

341              **Fig 6. SVC parameter selection result: (a) contour map (b) 3D view**

342    It can be seen from the results that the optimal values of $C$ and $g$ are 1.74 and 3.03, respectively. After

343    setting the parameters to the optimal values, we performed SVM modeling and obtained the predicted

344    values. Figure 7 shows the comparison charts of the actual and predicted values. After optimization, the

345    total number of support vectors is 19,460, and the number of support vectors at the boundary is 17,260.

346    After model training, the accuracy rate of the training set is 86.02%, and the accuracy rate of the test set

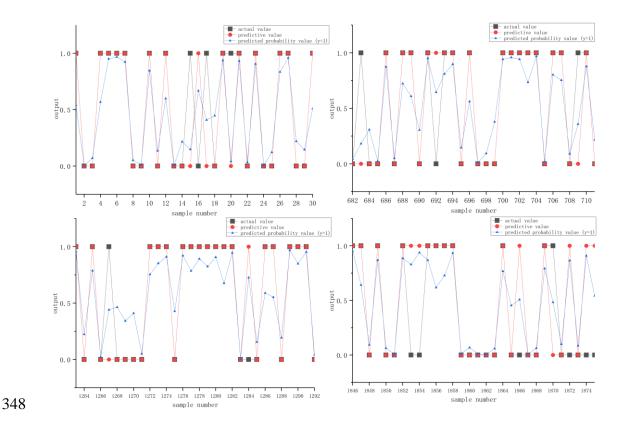347    is 84.27%, and the performance of the model is high.

348

**Fig. 7 Comparison charts of the predictive and actual values of the SVM (part of the**

**sample)**

### 3.2.4 Random Forest

We used the randomForest package in the R language to train random training samples. We then used

cross-validation to determine the optimal parameters of the model and the number of optimal decision

trees. Finally, we obtained the number of trees and the accuracy of the test and training data through

cross-validation. As shown in Figure 8, when the number of decision trees is 400, the accuracy tends to

be stable. We used the optimal number of decision trees to create the comparison charts of the actual

and predicted values of the test set (Figure 9) and the average accuracy decline of 13 forest fire driving

factors (Figure 10). It can be seen from Figure 10 that, among the main forest fire driving factors in

China, the location variables that have the greatest influence on the occurrence of forest fires are

longitude and latitude. Rainfall is the variable with the least influence on the occurrence of forest fires.
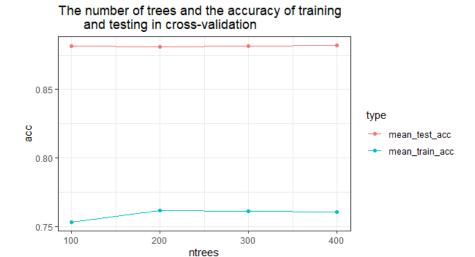
361

362    **Fig. 8 The number of trees and the accuracy of training and testing in cross-validation**
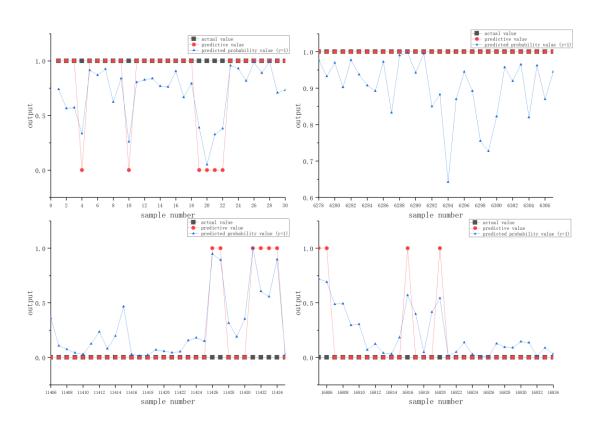


363

364    **Fig. 9 Comparison charts of the predictive and actual values of the RF (part of the sample)**
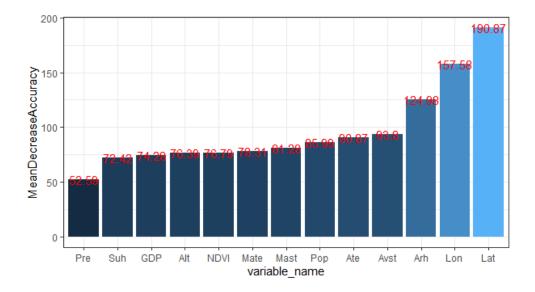
365

**Fig. 10 Mean decrease accuracy of 13 variables**

## 3.3 Accuracy Evaluation

We used the prediction results of the four models to construct a confusion matrix to obtain the accuracy,

precision, recall, f1 value, and AUC value, as shown in Table 6. Figure 11 shows the visualization of

the accuracy, precision, recall, and f1 values of the four models. Figure 12 shows the ROC curves of

the four models. The accuracy and f1 value of each model are more than 75%, and the AUC value is

more than 0.80. Thus, the performance of all four models is high. Among the four models, the RF

model has the highest predictive ability, with an accuracy rate of 89.2%, an f1 value of 89%, and the

highest AUC value, reaching 0.960. Compared with the other three models, the prediction ability of the

RBF neural network is the lowest, with an accuracy rate of 75.8% and an AUC value of 0.840. As

shown in Figures 11 and 12, the RF model outperforms the other three models. We therefore consider

the RF model to be the most suitable of the four models for forest fire prediction in China.

**Table 6: Evaluation results of the four models**

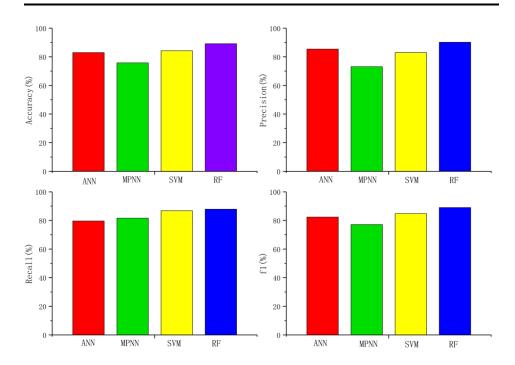| Model | Accuracy (%) | Precision (%) | Recall (%) | f1 value (%) | AUC |
|---|---|---|---|---|---|
| ANN | 83.0 | 85.4 | 79.6 | 82.4 | 0.904 |
| RBFN | 75.8 | 73.1 | 81.6 | 77.1 | 0.840 |
| SVM | 84.3 | 83.0 | 86.8 | 84.8 | 0.917 |
| RF | 89.2 | 90.2 | 87.9 | 89.0 | 0.960 |

379



380 **Fig 11. Comparison charts of accuracy, precision, recall, and f1 values of the four models**
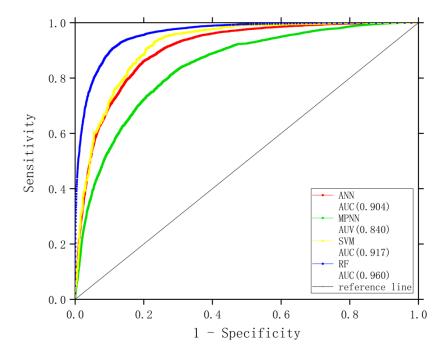
381

**Fig. 12 ROC curves of the four models**

## 3.4 Fire Risk Classification

384      After evaluating the accuracy of the four models, we used the RF model (highest accuracy) to obtain

385      the probability of forest fire occurrence for the full sample. We used ArcGIS to draw a forest fire

386      probability map (Figure 13) and a seasonal forest fire probability map (Figure 14) of China. Figure 13

387      shows that the high incidence of forest fires in China is mainly concentrated in the northeast (such as

388      the Greater Xing'an Mountains region), the southeast (such as Guangdong, Jiangxi, and Fujian), and

389      the southwest (such as Yunnan and Sichuan). On the whole, the probability of forest fires in eastern

390      China is higher than that in western regions, and the probability of forest fires in the north and south is

391      higher than that in central China. Figure 14 shows that the seasonal order of the probability of forest

392      fires in China is, from highest to lowest, spring, winter, summer, and autumn. Spring and winter are the

393      seasons with a high incidence of forest fires, which are mainly concentrated in northeast China (such as

394      Heilongjiang Province) and southeastern China (such as Fujian Province and Guangdong Province).

395



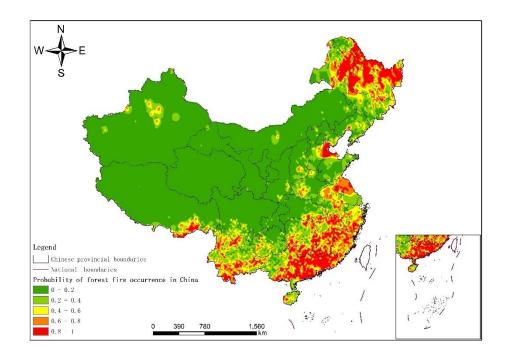396                    **Fig. 13 Forest fire probability map of China based on the RF model**
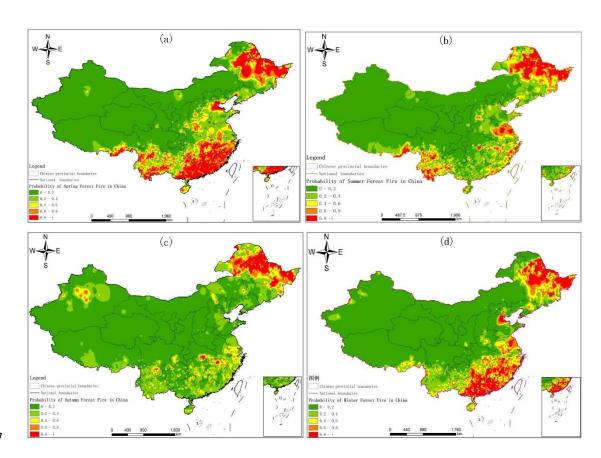
397



398          **Fig. 14 Seasonal forest fire probability map of China based on the RF model: (a) spring**

399    **(January, February, and March); (b) summer (April, May, and June); (c) autumn (July, August,**

400    **and September); and (d) winter (October, November, and December)**

## 4. Discussion

## 4.1 Major Forest Fire Driving Factors in China and Their Impact

403    In this study, we selected over 20 factors affecting the occurrence of forest fires. These factors can be

404    divided into six categories: location, meteorology, climate, terrain, society, and vegetation. Previous

405    research has been conducted on the impact of these factors on forest fires [64-67]. We used the RF

406    algorithm to filter the characteristics of these 20 forest fire driving factors and selected three or more

407    variables in the five variable subsets as the main forest fire driving factors. In this experiment there are

408    13 main forest fire driving factors: longitude, latitude, average surface temperature, daily maximum

409    surface temperature, cumulative precipitation, average relative humidity, sunshine hours, average

410    temperature, daily maximum temperature, altitude, population, GDP, and NDVI.

411    The results show that the location factors (longitude and latitude) are the factors that have the greatest

412    impact on the occurrence of forest fires. At present, there is little research that considers longitude and

413    latitude as the main driving factors for forest fires. Some researchers, however, have confirmed that, as

414    the latitude decreases, the number of forest fires increases [68-69]. China spans a large area, and the

415    two factors of longitude and latitude reflect the regional differences in forest fires in China. Forest fires

416    are more likely to occur at certain longitudes and latitudes.

417    Climate factors also have a great impact on forest fires, which is consistent with the findings of

418    previous studies [70-73]. Temperature is one of the three necessary conditions for combustion. When

419    the temperature reaches a certain level, forest fires are more likely to occur. The longer the sunshine

420    hours, the higher the temperature, and the greater the probability of forest fires. Rainfall and average

421    relative humidity are also among the main factors affecting forest fires [53][74-75].

422    Another type of driving factor is social and human factors (population and GDP). The larger the

423    population, the greater the human activity in the region and the more likely it is that human factors

424    cause forest fires. Catry et al. [76] (2007)and Sepulveda[77] (2001)reached the same conclusion.

425    Altitude and vegetation (NDVI) also affect the occurrence of forest fires to varying degrees. Tian et al.

426    [78] contend that forest fires mainly occur in low-altitude areas (2013). Chuvieco et al. [79] used NDVI

427    as a driving factor for forest fires to estimate fuel moisture (2004). The greater the NDVI value, the

428    higher the vegetation coverage; and the greater the flammability of the tree species, the more likely

429    such trees are to cause problems related to forest fires.

430    In this study, we used a RF algorithm for feature selection and eliminated seven variables: average

431    wind speed, maximum wind speed, aspect, slope, distance from fire point to highway, closest distance

432    from fire point to residential area, and special festivals. In this experiment, these factors have little

433    effect on the occurrence of forest fires. It may be that these factors change with time and space,

434    however. Other studies have found that these factors are among the main drivers of forest fires[80][81].

435    We believe that this may be due to the difference in the selected data and the difference in the method

436    of feature selection. In future research, a variety of feature-screening methods and analysis of different

437    regions may be used to obtain more comprehensive results.

438    ## 4.2 Optimal Choice of Forest Fire Prediction Model

439    We entered the forest fire driving factors selected by feature selection into the four models (ANN, RBF

440    neural network, SVM, and RF) for training. We then evaluated them using five criteria: accuracy,

441 precision, recall rate, f1 value, and AUC value. We selected the RF model as the optimal choice for

442 forest fire prediction. The accuracies of all four models are above 75%, which means that they are all

443 reliable. The RF model, however, exhibits the greatest prediction ability. The RBF neural network

444 model has the lowest prediction performance.

445 Samaher et al. [25] used a cascade correlation network, multilayer perceptron neural network,

446 polynomial neural network, RBF, and SVM for forest fire prediction(2018). They found the prediction

447 performance of the SVM was the highest, and the performance of RBF was the lowest, which is

448 consistent with our conclusion. Sakr et al. [36] used an SVM and an ANN to predict the fire risk in

449 Lebanon (2011). Their results showed that the performance of the SVM model was higher than that of

450 the ANN model. This finding is similar to ours. Paulo et al. [8] used RF and SVM models to predict

451 forest fires (2007). They concluded that the performance of the SVM model was higher than that of the

452 RF. Their finding is different from our results. A possible reason for this difference is that Paulo et al.

453 chose four types of weather factors and made predictions about small areas, whereas we chose 13 types

454 of factors and made predictions about a large area. The choice of variables and the difference in the

455 sample size affect the model training.

456 Bisquert et al. [82] used an ANN to establish a forest fire hazard model with a highest accuracy rate of

457 76%, which is lower than the accuracy of our model (83%) (2012). Hong [83] used an SVM algorithm

458 to analyze Dayu County in southwestern Jiangxi Province, China (2018). The results show that the

459 AUC value of the SVM is 0.75, which is lower than the value in our model (0.92). Pourtaghid et al.[45]

460 used an RF to create forest fire sensitivity analysis with a prediction accuracy of 72.8% (2016). Our

461 model reached a prediction accuracy of 89.2%.

462 The four models we selected all exhibit high predictive capabilities. The main reason may be that

463 appropriate multi-dimensional variables has been screened out and the data sample size is large, which

464 makes the training of each model more accurate and reliable.

465 There are also differences in the characteristics of these four models. The ANN and RBF neural

466 network models can be trained very quickly, and they can handle samples with a large amount of data,

467 but their accuracy in this experiment is relatively low. In subsequent studies, particle swarms or genetic

468 algorithms could be used to improve the accuracy of these models. The SVM model has a high

469 predictive ability, but it also has shortcomings. The higher the model complexity, the lower the

470 calculation speed. It takes a longer time in this model to obtain the optimal parameters when processing

471 large sample data. We will consider using other algorithms to optimize SVM in the future. The RF

472 model exhibited excellent expressive power in this experiment. It can quickly process large data

473 samples while ensuring high prediction accuracy.

474 ## 4.3 Recommendations for Forest Fire Prevention

475 We produced a probability map of forest fires in China that shows that the highest incidences of forest

476 fires are in the northeast (Heilongjiang Province and the northern Inner Mongolia Autonomous Region),

477 the southeast (Fujian Province, Guangdong Province, and Jiangxi Province), and Yunnan Province. The

478 pattern of forest fire points presents a spatial clustering distribution. Ma et al. obtained similar results

479 [4]. For these high-incidence areas, watch towers and monitoring equipment should be added for

480 monitoring and management. Moreover, the length of the forest fire barrier net should be increased to

481 reduce the spread of fires. In addition, the number of fire brigades and fire vehicles should be increased

482 to enhance the disaster-mitigation capabilities. Regarding seasonal risks of forest fires, forest fire

483   prevention and control should be emphasized in spring and winter. Strengthening fire-prevention

484   management during these periods would mainly involve strengthening the management of human

485   activities to reduce human-made forest fires and improving publicity and education, such as the

486   addition of fire-prevention signs.

487   This study has some shortcomings, and there is scope for improvement. One of the three elements of

488   fire is combustible fuel. For the selection of forest fire driving factors, however, there is currently no

489   way to obtain data on fuel load and other related factors. Thus, this experiment lacked relevant data

490   such as combustible load, particle size of combustible material, and combustible tree species. If

491   possible, in future research, such data could be added to the forest fire prediction model.

492   This study selected four kinds of machine-learning algorithms for the forest fire prediction model.

493   Other applicable machine-learning algorithms could be used in future experiments. In addition, the

494   ability of these machine-learning algorithms to analyze spatial heterogeneity is relatively weak.

495   Subsequent research could use geographically weighted regression to build a high-precision forest fire

496   prediction model.

## 497   5. Conclusion

498   This study determined the main driving factors of forest fire occurrence in China through feature

499   selection. The main factors that affect the occurrence of forest fires to varying degrees are

500   meteorological, topographic, man-made, and vegetation factors. We built four forest fire prediction

501   models using the following machine-learning algorithms: an artificial neural network, a radial basis

502   function neural network, a support-vector machine, and a random forest. The results of the evaluation

503   show that the accuracy of all the models is higher than 75%. These models can be used to build forest

504  fire prediction models. Among the four models, the RF model has the highest comprehensive predictive

505  ability, with an accuracy of 89.25%. It is therefore the optimal choice for a forest fire prediction model

506  in China.

507  We used the RF model to predict the probabilities of forest fires in China. Based on these probabilities,

508  we drew a map of the probability of forest fire occurrence in China and a map of the probability of

509  forest fires in China by season (spring, summer, autumn, and winter). Finally, based on these maps, we

510  identified the high-incidence areas and areas at risk of forest fires. We then put forward fire-prevention

511  recommendations for the corresponding regions and seasons.

512  This research helps in understanding the main forest fire driving factors in China. It provides a

513  reference for the selection of high-precision forest fire prediction models. In addition, it provides

514  suggestions on the time and location of forest fire prevention in China. Moreover, this study provides

515  guidance for China's forest fire prevention and control work.

516  ## Data Availability

517  The data used to support the findings of this study are available from the corresponding author upon

518  request.

519  ## Conflicts of Interest

520  The authors declare that they have no conflicts of interest.

## Funding Statement

## Authors' contributions

Yudong Li performed the experiment and wrote the manuscript;

Zhongke Feng contributed to the conception of the study;

Ziyu Zhao and Shilin Chen helped perform the analysis with constructive discussions;

Hanyue Zhang helped perform the data analysis.

## Acknowledgments

## References

[1]  Bhusal, Satish & Mandal, Ram. (2020) Forest fire occurrence, distribution and future risks in Arghakhanchi district, Nepal. Journal of Geography. 2(1):10-20.

[2]  B. K. Singh, N. Kumar and P. Tiwari. (2019) "Extreme Learning Machine Approach for Prediction of Forest Fires using Topographical and Metrological Data of Vietnam," 2019 Women Institute of Technology Conference on Electrical and Computer Engineering (WITCON ECE), Dehradun Uttarakhand, India,. pp. 104-112. doi: 10.1109/WITCONECE48374.2019.9092926.

[3]  D Mandallaz, R Ye. (1997) Prediction of forest fires with Poisson models, Canadian Journal of

541     Forest Research. 27(10): 1685-1694.https://doi.org/10.1139/x97-103.

542     [4]  Ma, W., Feng, Z., Cheng, Z., Chen, S., & Wang, F. (2020) Identifying Forest Fire Driving Factors

543          and Related Impacts in China Using Random Forest Algorithm. Forests. 11(5), 507.

544          doi:10.3390/f11050507.

545     [5]  Dimopoulou, M.; Giannikos, I. (2004) Towards an integrated framework for forest fire control.

546          Eur. J. Oper. Res. 152: 476–486.

547     [6]  Flannigan, M.D.; Krawchuk, M.A.; Groot, W.J.D.; Wotton, B.M.; Gowman, L.M. (2009)

548          Implications of changing climate for global wildland fire. Int. J. Wildland Fire. 18:483–507.

549     [7]  Z. Q. Xu, X. Y. Su, Y. Zhang. (2012) Forest Fire Prediction Based on Support Vector Machine.

550          Chinese Agricultural Science Bulletin.28(13):126-131.

551     [8]  Paulo Cortez, Anibal Morais. (2007) A Data Mining Approach to Predict Forest Fires using

552          Meteorological Data. New Trends in Artificial Intelligence. 512-523.

553     [9]  Barik A, Baidya Roy S. (2020) Effects of meteorology on forest fires in India: A modeling

554          study[C]//EGU General Assembly Conference Abstracts. 18317.

555     [10] Chetehouna K. (2020) Overview of The Forest Fire Research[C]//Proceeding International

556          Conference on Science and Engineering. 3: xvii-xvii.

557     [11] Shang C, Wulder M A, Coops N C, et al. (2020) Spatially-Explicit Prediction of Wildfire Burn

558          Probability Using Remotely-Sensed and Ancillary Data. Canadian Journal of Remote Sensing.

559          1-17.

560     [12] Z.W. SU, A.Q. LIU et al. (2016) Driving factors and spatial distribution patteren of forest fire in

561          Fujian Province. JOURNAL OF NATURAL DISASTERS. 25(2):110-119.

562     [13] M. Yu, (2016) The establishment and study of Regional forest fire prediction and forecasting

563    model. Beijing, Beijing Forestry University.

564    [14] Y. Song. (2018) Research on Forest Fire Drivers and Models in Heilongjiang Province. Herbing,

565    Northeast Forestry University.

566    [15] CV Garcia, PM Woodard, SJ Titus, WL Adamowicz and BS Lee. (1995) A Logit Model for

567    Predicting the Daily Occurrence of Human Caused Forest-Fires. International Journal of Wildland

568    Fire. 5(2):101-111.

569    [16] H. Peng , M.C. Shi , Y. Sun, Mike Wotton. (2014) Lightning Fire Forecasting Model of Daxing'an

570    Mountain Based on Logistic Model. Journal of Northeast Forestry University. 42(7):166-169.

571    [17] W.G. Wang, J.H. Pan, Y.Y Feng, Z. Li, L. L. Dong. (2017) Model and Zoning of Fire Risk in

572    Gansu Province based on GWLR and MODIS Imagery. Remote Sensing Technology and

573    Application. 32(03):514-523.

574    [18] Futao Guo, Selvara Selvalakshmi, Fangfang Lin, Guangyu Wang, Wenhui Wang, Zhangwen Su,

575    Aiqin Liua.(2016)Geospatial information on geographical and human factors improved

576    anthropogenic fire occurrence modeling in the Chinese boreal forest. Canadian Journal of Forest

577    Research. 46(4): 582-594. https://doi.org/10.1139/cjfr-2015-0373.

578    [19] K.Z. Liu, L.F. Shu, F.J. Zhao et al. (2017) Research on spatial distribution of forest fire based on

579    satellite hotspots data and forecasting model. Journal of Forestry Engineering, 2:128-133.

580    [20] B.Q. LIAO, J. WEI et al. (2008) Logistic and ZIP Regression Model for Forest Fire Data. FIRE

581    SAFETY SCIENCE.3:143-149.

582    [21] D Mandallaz, R Ye. (1997) Prediction of forest fires with Poisson models. Canadian Journal of

583    Forest Research.27(10): 1685-1694. https://doi.org/10.1139/x97-103.

584    [22] F.T. GUO, H.Q. HU et al. (2010) Relationship between forest lighting fire occurrence and weather

585    factors in Daxing'an Mountains based on negative binomial model and zero-inflated negative

586    binomial models. Chinese Journal of Plant Ecology.21(01):159-164.

587    [23] H.L. Liang ,Y.R. Lin, G. Yang et al. (2016) Application of Random Forest Algorithm on the Forest

588    Fire Prediction in Tahe Area Based on Meteorological Factors. SCIENTIA SILVAE

589    SINICAE.52(01):89-98.

590    [24] Z.Q. Xu. (2012) Forest Fire Area Prediction Based on Support Vector Machine. Beijing, Beijing

591    Forestry University.

592    [25] Samaher Al_Janabia et al. (2018) Assessing the suitability of soft computing approaches for forest

593    fires prediction. Applied Computing and Informatics. (14): 214-224.

594    [26] Volkan Sevinca, Omer Kucukb, Merih Goltasc. (2020) A Bayesian network model for prediction

595    and analysis of possible forest fire causes. Forest Ecology and Management. (457):1-11.

596    [27] T. Artés, A. Cencerrado, A. Cortés, T. Margalef. (2016) Time aware genetic algorithm for forest

597    fire propagation prediction: exploiting multi-core platforms, Concurrency and Computation:

598    Practice and Experience.

599    [28] D.T. Bui, Q.T. Bui, Q.P. Nguyen, B. Pradhan, H. Nampak, P.T. Trinh, (2017) A hybrid artificial

600    intelligence approach using GIS-based neural-fuzzy inference system and particle swarm

601    optimization for forest fire susceptibility modeling at a tropical area, Agric. For. Meteorol. 233:

602    32-44

603    [29] M. Denham, A. Cortés, T. Margalef, E. Luque, (2008) Applying a dynamic data driven genetic

604    algorithm to improve forest fire spread prediction. International Conference on Computational

605    Science. Springer Berlin Heidelberg. pp.36-45.

606    [30] H. Hong, S.A. Naghibi, M.M. Dashtpagerdi, H.R. Pourghasemi, W. Chen, (2017) A comparative

607  assessment between linear and quadratic discriminant analyses (LDA-QDA) with frequency ratio

608  and weights-of-evidence models for forest fire susceptibility mapping in China. Arab. J. Geosci.10

609  (7): 167.

610 [31] A.M. Özbayog ̆lu, R. Bozer, (2012) Estimation of the burned area in forest fires using

611  computational intelligence techniques, Proc. Comput. Sci. 12 :282–287.

612 [32] H. Soliman, K. Sudan and A. Mishra. (2010) A smart forest-fire early detection sensory system:

613  Another approach of utilizing wireless sensor and neural networks. SENSORS, 2010 IEEE, Kona,

614  HI. pp. 1900-1904.doi: 10.1109/ICSENS.2010.5690033.

615 [33] Çetin Elmas, Yusuf Sönmez. (2011) A data fusion framework with novel hybrid algorithm for

616  multi-agent Decision Support System for Forest Fir. Expert Systems with

617  Applications.38(8):9225-9236. https://doi.org/10.1016/j.eswa.2011.01.125.

618 [34] Onur Satir, Suha Berberoglu & Cenk Donmez (2016) Mapping regional forest fire probability

619  using artificial neural network model in a Mediterranean forest ecosystem, Geomatics. Natural

620  Hazards and Risk. 7: 1645-1658, doi: 10.1080/19475705.2015.1084541.

621 [35] Maeda, E. E., Formaggio, A. R., Shimabukuro, Y. E., Arcoverde, G. F. B., & Hansen, M. C.

622  (2009). Predicting forest fire in the Brazilian Amazon using MODIS imagery and artificial neural

623  networks. International Journal of Applied Earth Observation and Geoinformation, 11(4), 265–

624  272. doi: 10.1016/j.jag.2009.03.003.

625 [36] Sakr G E, Elhajj I H, Mitri G. (2011) Efficient forest fire occurrence prediction for developing

626  countries using two weather parameters. Engineering Applications of ArtificialIntelligence.24(5):

627  888-894.

628 [37] B.C. Ko, K.H. Cheong, J.Y. Nam. (2009) Fire detection based on vision sensor and support vector

629       machines. Fire Saf. J. 44 (3): 322-329.

630    [38] G.E. Sakr, I.H. Elhajj, G. Mitri, Efficient forest fire occurrence prediction for developing countries

631       using two weather parameters. Eng. Appl. Artif. Intell. 24.

632    [39] D.W. Xie, S.L. Shi,(2014)Prediction for burned area of forest fires based on SVM model, in:

633       Applied Mechanics and Materials. Trans Tech Publications. 513(5):4084-4089.

634    [40] J. Zhao, Z. Zhang, S. Han, C. Qu, Z. Yuan, D. Zhang. (2011) SVM based forest fire detection

635       using static and dynamic features. Computer Sci. Inform. Syst. 8(3): 821–841.

636    [41] Cutler DR , et al. (2007) Random forests for classification in Ecology. Ecology.88(11):2783-2792.

637    [42] Prasad AM, et al.  (2006) Newer classification and regression tree techniques: Bagging and

638       random forests for ecological prediction. Ecosystems .9(2):181-199.

639    [43] Rodrigucs M, Dc la Riva J.  (2014) An insight into machines learning algorithms to model

640       humarrcaused wildfire or currence. Environmental Modelling & Software,57:192-201.

641    [44] LIANG H L, LIN Y R, YANG G, et al.(2016) Application of random forest algorithm on the

642       forest fire prediction in Tahe area based on meteorological factors. Scientla Silvae

643       Sinicae.52(1):89-98.

644    [45] Pourtaghi, Z. S., Pourghasemi, H. R., Aretano, R., & Semeraro, T. (2016) Investigation of general

645       indicators influencing on forest fire and its susceptibility modeling using different data mining

646       techniques. Ecological Indicators. 64:72–84. doi:10.1016/j.ecolind.2015.12.030.

647    [46] Tian, X., Zhao, F., Shu. L. (2013) Wang, M. Distribution characteristics and the influence factors

648       of forest fires in China. For. Ecol. Manag. 310:460–467.

649    [47] Chang, Y., Zhu, Z., Bu, R., Li, Y., Hu, Y. (2015) Environmental controls on the characteristics of

650       mean number of forest fires and mean forest area burned (1987–2007) in China. For. Ecol. Manag.

651      356:13–21.

652      [48] Zhong. M., Fan, W., Liu. T., Li. P. (2003) Statistical analysis on current status of China forest fire

653      safety. Fire Saf. J. 38:257–269.

654      [49] Aifeng, L.U. (2011) Study on the relationship among forest fire, temperature and precipitation and

655      its spatial-temporal variability in China. Agric. Sci. Technol. 12:1396–1400.

656      [50] Ying. L., Han. J., Du. Y., Shen, Z. (2018) Forest fire characteristics in China: Spatial patterns and

657      determinants with thresholds. For. Ecol. Manag.424:345–354.

658      [51] X.K. Ren, Chinese Geography, Chinese Press, ISBN:9787507528800.

659      [52] F.T. Guo, Z.W. Su, G.Y. Wang et al. (2017) Understanding fire drivers and relative impacts in

660      different Chinese forest ecosystems. Science of the Total Environment.605: 411-425.

661      [53] GUO F, WANG G, SU Z, et al. (2016) What drives forest fire in Fujian, China? Evidence from

662      logistic regression and random forests. International Journal of Wildland Fire.25(5):505-519.

663      [54] Chang Y, Zhu Z L, Bu R C, et al. (2013) Predicting fire occurrence pat-terns with logistic

664      regression in Heilongjiang Province, China. Landscape Ecology.28(10) : 1989 -2004 .

665      [55] XU X L. (2018) Spatial distribution data set of quarterly vegetation index (NDVI) in China. Data

666      Registration and Publishing System of Resources and Environmental Science Data Center of

667      Chinese Academy of Sciences (http://www.resdc.cn/DOI).

668      [56] W.Y MA, Z.K. FENG et al. (2020) Study on driving factors and distribution pattern of forest fires

669      in Shanxi province. Journal of Central South University of Forestry & Technology.1-13.

670      [57] ERIDUNTONGLAGA. (2013) Analysis on landscape pattern of woodland based on DEM in the

671      Daqingshan Mountains Inner Mongolia. Huhhot: Inner Mongolia Normal University.

672      [58] A. Subasi, (2007) EEG signal classification using wavelet feature extraction and a mixture of

673    expert model. Expert Syst. Appl. 32 (4): 1084-1093.

674    [59] D. Chen. (2019) Prediction of Forest Fire Occurrence in Daxing'an Mountains Based on Logistic

675    Regression Model," FOREST RESOURCES MANAGEMENT. (01):116-122.

676    [60] Q.Q. Wang, J.T. Tong, L. Zhang et al. (2020) Seismic data denoising using multi-layer perceptron.

677    Oil Geophysical Prospecting.55(02):272-281+228.

678    [61] D.S. Broomhead, D. Lowe, (1988) Radial basis functions, multi-variable functional interpolation

679    and adaptive networks (No. RSRE-MEMO-4148). Royal signals and radar establishment Malvern

680    (United Kingdom).

681    [62] X.C. Wang et al. (2013) 43 Cases of MATLAB neural network analysis. Beijing, Bei hang

682    University Press.

683    [63] Yudong Li, Zhongke Feng , Shilin Chen, Ziyu Zhao, and Fengge Wang. (2020) Application of the

684    Artificial Neural Network and Support Vector Machines in Forest Fire Prediction in the Guangxi

685    Autonomous Region, China. Discrete Dynamics in Nature and Society.

686    https://doi.org/10.1155/2020/5612650.

687    [64] Ganteaume A, Camia A, Jappiot M, et al. (2013) A review of the main driving factors of forest fire

688    ignition over Europe. Environmental Management. 51(3): 651-662.

689    [65] Syphard AD, Radeloff VC, Keuler NS, Taylor RS, Hawbaker TJ, Stewart SI, Clayton MK (2008)

690    Predicting spatial patterns of fire on a southern California landscape. Int J Wildland Fire 17:602–

691    613.

692    [66] Pew KL, Larsen CPS (2001) GIS analysis of spatial and temporal patterns of human-caused

693    wildfires in the temperate rainforest of Vancouver Island, Canada. Forest Ecol Manag 140:1-18.

694    [67] Dickson BG, Prather JW, Xu Y, Hampton HM, Aumack EN, Sisk TD (2006) Mapping the

695    probability of large fire occurrence in northern Arizona. USA Landsc Ecol 2:747-761.

696    [68] Ying. L., Han. J.; Du. Y., Shen, Z. (2018) Forest fire characteristics in China: Spatial patterns and

697    determinants with thresholds. For. Ecol. Manag. 424:345–354.

698    [69] Prasad, A.M.; Iverson, L.R.; Liaw, A. (2006) Newer classification and regression tree techniques:

699    Bagging and random forests for ecological prediction. Ecosystems. 9:181-199.

700    [70] Guo F, Su Z, Wang G, et al. (2016) Wildfire ignition in the forests of southeast China: Identifying

701    drivers and spatial distribution to predict wildfire likelihood. Applied Geography. 12-21.

702    [71] Liu, Z., Yang, J., Chang, Y., Weisberg, P. J., & He, H. S. (2012). Spatial patterns and drivers of fire

703    occurrence and its future trend under climate change in a boreal forest of northeast China. Global

704    Change Biology, 18:2041-2056.

705    [72] Syphard, A. D., Radeloff, V. C., Keuler, N. S., Taylor, R. S., Hawbaker, T. J., Stewart, S. I., et al.

706    (2008) Predicting spatial patterns of fire on a southern California landscape. International Journal

707    of Wildland Fire.17:602-613.

708    [73] Ying L, Han J, Du Y, et al. (2018) Forest fire characteristics in China: Spatial patterns and

709    determinants with thresholds. Forest Ecology and Management.345-354.

710    [74] Zumbrunnen, T., Pezzatti, G. B., Men?endezd, P., Bugmann, H., Bürgi, M., & Conedera, M. (2011)

711    Weather and human impacts on forest fires: 100 years of fire history in two climatic regions of

712    Switzerland. Forest Ecology and Management.261:2188-2199.

713    [75] Cardille J A, Ventura S J, Turner M G, et al. (2001) ENVIRONMENTAL AND SOCIAL

714    FACTORS INFLUENCING WILDFIRES IN THE UPPER MIDWEST, UNITED STATES.

715    Ecological Applications.11(1): 111-127.

716    [76] Catry F X, Damasceno P, Silva J S. (2007) Spatial distribution patterns of wildfire ignitions in

717    Portugal. Available at: http://www. eufirelab.org/toolbox2/library/upload/2380.pdf, 2009-08.

718    [77] Sepulveda B J I, Meza S R, Zuniga C W R et al. (2001) GIS to Determine Risk of Forest Fires in

719    Northwestern Mexico. Mexico: Technical Publication. 37.

720    [78] Tian, X.; Zhao, F.; Shu, L.; Wang, M. (2013) Distribution characteristics and the influence factors

721    of forest fires in China. For. Ecol. Manag. 310: 460-467.

722    [79] Chuvieco, E.; Cocero, D.; Riaño, D.; Martin, P .; Martínez-Vega, J.; de la Riva, J.; Pérez, F. (2004)

723    Combining ndvi and surface temperature for the estimation of live fuel moisture content in forest

724    fire danger rating. Remote Sens. Environ. 92: 322-331.

725    [80] Maingi K J, Henry M C. (2007) Factors influencing wildfire occurrence and distribution in eastern

726    Kentucky, USA. International Journal of Wildland Fire, 16(1): 23-33. doi: 10.1071/ WF06007.

727    [81] Avilaflores D Y, Pompagarcia M, Antonionemiga X, et al. (2010) Driving factors for forest fire

728    occurrence in Durango State of Mexico: A geospatial perspective. Chinese Geographical Science.

729    20(6): 491-497.

730    [82] Bisquert M, Caselles E, Sanchez J M, et al. (2012) Application of artificial neural networks and

731    logistic regression to the prediction of forest fire danger in Galicia using MODIS data.

732    International Journal of Wildland Fire. 21(8): 1025-1029.

733    [83] Hong, H., et al., (2018) Applying genetic algorithms to set the optimal combination of forest fire

734    related variables and model forest fire susceptibility based on data mining models. The case of

735    Dayu County, China. Science of The Total Environment. 630:1044-1056.